



THE UNIVERSITY *of* EDINBURGH

## Edinburgh Research Explorer

### Video Anomaly Detection in Real-Time on a Power-Aware Heterogeneous Platform

**Citation for published version:**

Blair, C & Robertson, N 2015, 'Video Anomaly Detection in Real-Time on a Power-Aware Heterogeneous Platform', *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 26, no. 11, pp. 2109-2122. <https://doi.org/10.1109/TCSVT.2015.2492838>

**Digital Object Identifier (DOI):**

[10.1109/TCSVT.2015.2492838](https://doi.org/10.1109/TCSVT.2015.2492838)

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Peer reviewed version

**Published In:**

IEEE Transactions on Circuits and Systems for Video Technology

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



# Video Anomaly Detection in Real-Time on a Power-Aware Heterogeneous Platform

Calum G Blair and Neil M Robertson, Senior Member, IEEE<sup>‡</sup>

October 21, 2015

## Abstract

FPGAs and GPUs are often used when real-time performance in video processing is required. An accelerated processor is chosen based on task-specific priorities (power consumption, processing time and detection accuracy), and this decision is normally made once at design time. All three characteristics are important, particularly in battery-powered systems. Here we propose a method for moving selection of processing platform from a single design-time choice to a continuous run-time one. We implement Histogram of Oriented Gradients (HOG) detectors for cars and people and Mixture of Gaussians (MoG) motion detectors running across FPGA, GPU and CPU in a heterogeneous system. We use this to detect illegally parked vehicles in urban scenes. Power, time and accuracy information for each detector is characterised. An anomaly measure is assigned to each detected object based on its trajectory and location, when compared to learned contextual movement patterns. This drives processor and implementation selection, so that scenes with high behavioural anomalies are processed with faster but more power-hungry implementations, but routine or static time periods are processed with power-optimised, less

---

<sup>\*</sup>C. Blair is with the Institute for Digital Communications, University of Edinburgh, UK.

e-mail: c.blair@ed.ac.uk.

<sup>†</sup>N. Robertson is with the Visionlab, Heriot-Watt University.

<sup>‡</sup>Manuscript received xxx.

accurate, slower versions. Real-time performance is evaluated on video datasets including i-LIDS. Compared to power-optimised static selection, automatic dynamic implementation mapping is 10% more accurate but draws 12W extra power in our testbed desktop system.

## 1 Introduction

Surveillance systems have recently become more commonplace, more portable, and capable of more complex tasks. The process of selecting or designing a processing system for any surveillance application involves meeting some performance constraints (computation time must be fast enough for real-time performance) and optimising others (power consumption should be minimised and algorithm accuracy maximised). Improving any one of these characteristics will result in tradeoffs in the others. Here, we propose a real-time video surveillance system for detecting anomalous behaviour in an urban setting, and analyse its performance. In this scenario our video analysis system contains a heterogeneous set of processors and its power supply will be constrained in some way, such as relying on batteries. The ability to perform timely and accurate detections with minimised power consumption is important. At different time periods, power should be prioritised more than speed (or framerate), and vice versa, depending on the level of any perceived threat. This allows a power-constrained, real-time system to react quickly to relevant changes in its environment, while conserving power in periods of inactivity.

### 1.1 Motivation

In recent years, algorithms of increasing computational complexity have been developed, allowing more accurate detection or classification of objects of interest. In conjunction with the rise in pixel data volume from higher-resolution imaging devices, this has led to requirements to process higher volumes of data in a timely fashion — ideally in real-time. Many of these algorithms are ‘embar-

massively parallel', so two technologies have been employed to allow accelerated processing of video data, for online and offline real-time tasks [1, 2]. These devices represent two different approaches to the problem of parallelisation; FPGAs (Field-Programmable Gate Arrays) allow temporal parallelisation by building up long pipelines of simple arithmetic operations specifically matched to the problem in question, and permit spatial parallelisation by duplicating these units as broadly as possible to process multiple parts of the image at once. GPUs (Graphics Processing Units) use thousands of existing cores optimised for fast arithmetic to massively parallelise the calculations required by an algorithm. These approaches have tradeoffs; FPGAs have lower power consumption, but require considerably more time and knowledge to program [3]. They also have a long history and acceptance of being deployed in military systems. GPUs are programmed using a more familiar software design flow and can allow greater numerical precision due to their native reliance on floating-point arithmetic, but draw considerably more power.

The choice of platform to deploy an algorithm onto is typically made at design time, as part of a decision about how a task should be partitioned onto hardware. Typically, once a decision to deploy to a specific platform is made, the consequences (higher power consumption, reduced detection capability, etc.) remain part of that system and must be accepted throughout its use. Delaying the choice of processing architecture until the system is deployed allows the system to trade off power consumption and frame processing time *on-the-fly, in response to changes in scene conditions*. For example, when no activity is seen, or no anomaly is present, we can conserve power at the expense of frame rate and processing time. If a higher level of anomaly is present (for example a pedestrian entering an unusual area, or a car driving the wrong way down a road or parking in a forbidden location), any of these anomalous actions could represent a threat. For vehicle-mounted systems performing a surveillance task, this may mean a threat to its occupants, so in this case we would ideally process data and identify anomalies more quickly; immediate power consumption is less



important. These priorities change dynamically with scene content.

## 1.2 Contributions

In this paper we describe and analyse the performance of our heterogeneous system which uses CPU, FPGA and GPU to allow real-time detection of objects and anomalies. It autonomously chooses which set of processors to run an algorithm or algorithm stage on. This decision is made dynamically and depends on the level of anomaly seen in the video stream. We evaluate this on two video datasets and demonstrate detection of this behaviour in both datasets while reducing overall power consumption. Vehicles which park or stop in forbidden or unusual areas are used to represent anomalous behaviour. An example is shown in Fig. 1b, where vehicles are forbidden from parking on the double yellow lines at the roadside. An anomalous vehicle is denoted by a red box.

Our contributions are as follows:

1. We describe a hardware platform of heterogeneous processors (FPGA, GPU and CPU), with multiple possible implementations of computationally expensive algorithms which can be run on it. We characterise each implementation's power consumption, processing time and detection accuracy. Our implementations and their platforms can be switched between dynamically.
2. We introduce an algorithm for quantising the level of behavioural anomaly in a video stream using object and motion detectors.
3. We use a hardware mapping algorithm which uses anomaly level to prioritise reduced power consumption or shorter processing time. This then dynamically selects algorithm implementations using their performance characteristics. This set of implementations is used to process subsequent video data.

These contributions are shown in the overall system diagram in Fig. 2.

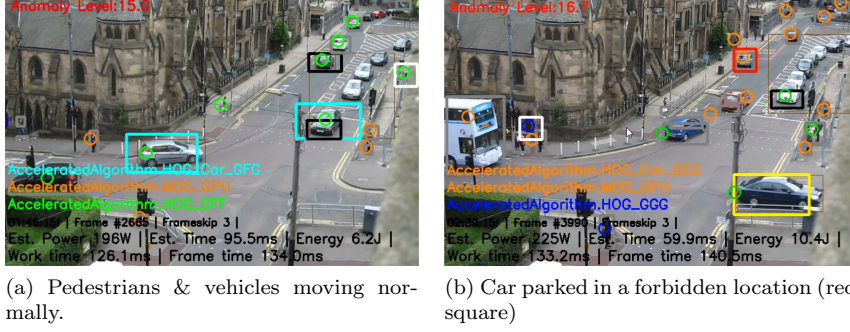


Figure 1: Anomaly detection in an urban scene (BankSt dataset). When an anomalous event (car stopping in an unusual area) is detected then detection speed is prioritised and power consumption increases as a result.

### 1.3 Paper Layout

This paper is structured as follows: Section 2 places this work in context by describing related work in the areas of anomaly detection, parallelised algorithm implementations and power-aware computing. Section 3 describes our system at a high level, and details the components used to perform object and anomaly detection. It also covers our mechanism for dynamically switching between algorithm implementations. Section 4 describes the video datasets and evaluation procedure and presents results. Section 5 discusses these in more detail, along with our system’s limitations. Finally, Section 6 restates our results and describes directions for future work.

## 2 Related Work

This paper is concerned with algorithms at multiple levels of abstraction, from feature extraction in images up to behaviour analysis. We aim to process video data in real time so here we also describe accelerated processing and how to schedule or allocate processing tasks to minimise power consumption.

## 2.1 Anomaly Detection

Loy *et al.* define anomalous events or unusual behaviour in video as being in one of three categories (those very different from the training set, those which are ambiguous and rarely present in the training set, and those with only weak visual evidence). Humans can fail to detect the latter two types, especially during longer tasks [4]. A review paper by Morris and Trivedi [5] considers surveillance in urban scenes and notes that as scenes become more unstructured, identifying unusual events becomes harder. Detection in urban scenes, with multiple classes of traffic participants and fewer restrictions on directions, is more difficult than behaviour recognition on a highly structured motorway, particularly when vehicles in urban environments may be stationary for long time periods. Anomaly detection spans several levels of abstraction, from individual object detections which can be grouped into trajectories, then clustered into common paths in a scene [6]. This produces clear trajectories but requires offline analysis. Piciarelli and Foresti describe a method for online clustering of object locations into trajectories [7]. A sliding temporal window allows tracks to be matched to clusters; the window expands as track length increases. Given a learned set of tracks clustered into trajectory trees, the likely behaviour of a new object can be predicted, as the most common clusters and parent-child cluster transitions are known [7]. Morris and Trivedi describe a system for detecting U-turn events in traffic [8], but we look here at the problem of detecting parked vehicles.

Albiol *et al.* [9] use spatiotemporal maps to identify these events in the i-LIDS dataset. This relies on human operators labelling regions where parking is forbidden. Regions are evaluated as belonging to the background or not, and so distinct objects are not detected or classified. Lee *et al.* also process some events from i-LIDS, using a 1-D transformation to detect vehicles parking illegally, and demonstrate real-time performance on CPU [10].

## 2.2 Accelerated Object Detection

Work on real-time object detection has attracted much interest recently, particularly given the possibilities for detection from embedded systems fitted to vehicles. A common pedestrian detection algorithm is HOG (Histogram of Oriented Gradients) [11]. HOG is split into three computationally expensive algorithm stages: image resizing, feature extraction and classification. Many improvements and variations to the overall algorithm are documented in a review paper by Dollar *et al.* [12]. Recent work has demonstrated that performing feature extraction on fewer scales and rescaling the resulting features is also effective [13]. Current state-of-the-art detectors are still based on this work; they extract various colour and shape features before classifying these using a Support Vector Machine (SVM) or boosted classifier. Multiple hardware implementations of HOG and its derivatives have been proposed, on various devices including GPU [1, 14], FPGA [15] and a hybrid system [16], with feature extraction done on the FPGA and SVM classification handled by the GPU. A comparison of these platforms in a heterogeneous system is performed by Blair *et al.* [17], where processing time, detection accuracy and system power consumption of the HOG algorithm are all compared. The image resize, feature extraction and SVM classification subtasks can be allocated to GPU, FPGA or CPU. The algorithm runs fastest when all tasks are performed on GPU, and uses the least power when CPU and FPGA are used. A mix of GPU and FPGA processing allows a reasonable tradeoff between power and framerate.

A car detection algorithm has been built using a similar approach, using HOG with deformable part models [18].

## 2.3 Platform Selection and Power-aware Computing

FPGAs and GPUs occupy different points in design space. Allocating tasks (those defined as stages of an algorithm or algorithms in their entirety) onto processors is a form of *design space exploration*. Work by Cope *et al.* has com-

pared FPGA and GPU for accelerating low-level image processing operations such as convolutions and enhancement algorithms [19]. GPU performance decreases with increasing kernel size, but there is no single platform which is always the most suitable and the data flow of the desired algorithm must always be considered. Recent work by Wu *et al.* reinforces these conclusions, while taking power into account [20].

Points in design space are evaluated for optimality using a Pareto curve [21]. Optimal points dominate non-optimal ones if they improve on at least one characteristic. Eventually a Pareto front of all the optimal points is formed. Prioritising one characteristic over another will move between these points.

Yu and Prasanna [22] explore power-aware resource allocation for static tasks. Quinn *et al.* explored tradeoffs in assigning discrete algorithm stages to FPGA or CPU for pixel processing algorithms [23]. Tradeoffs were expressed in terms of chip area versus instruction latency, rather than accuracy and power. However, there is a limited body of work on using power consumption information to map tasks to processors in video. The main work we are aware of in this field is by Llamocca and Pattichis [24]. They perform dynamic multiobjective optimisation in a power-performance-accuracy (PPA) space (in this context, performance refers to processing time). This is done for pixel processing algorithms including gamma correction and contrast enhancement, where the design space is populated by implementations with different clock frequencies and bit widths. A user-chosen accuracy level drives dynamic selection, and the selected implementation is loaded via FPGA dynamic reconfiguration.

A conference version of this paper [25] considered power-aware dynamic allocation on algorithms operating at a higher level of abstraction, where multiple heterogeneous accelerators are used. We expand on this work here by adding another video evaluation, further algorithm and implementation details and additional experimental results, including time taken to switch between selected implementation mappings.

Table 1: Algorithms and Implementations Used.

Algorithm (* means compute-intensive)	Implementation(s)		
Pedestrian Detection*: PED-HOG (§3.2)	FPGA	GPU	CPU
Car Detection*: CAR-HOG (§3.3)	FPGA	GPU	CPU
Motion Detection*: MoG (§3.4)		GPU	
Tracking: Kalman Filter (§3.5)			CPU
Trajectory Clustering (§3.6)			CPU
Bayesian Motion Context (§3.7)			CPU

### 3 System Details

Our system detects objects and motion in a video view of a scene, uses contextual information obtained from previous behaviour to assign an anomaly level to scene activity, and uses this to select the optimal algorithm implementations to process the next incoming frame. We detect both pedestrians and cars; this allows identification of common traffic participants and hence their behaviours within the environment. The system processes offline videos, but operates at 25 frames per second (fps) and dynamically calculates the number of frames to drop to maintain this realtime processing rate. Time to decode the input video and time to display or annotate an output image are not counted within this window (as we assume that live video would be captured without needing to be decoded, and the only outputs required are infrequent snapshots of anomalous events as they occur).

Fig. 2 gives a high-level overview of our system. Where relevant, we present performance information on aspects of the system in this section; overall results for the anomaly detection task are presented in Section 4.

In Fig. 2, the ‘detection algorithms’ block processes each image and generates bounding boxes describing object location and classification. These are used by the later algorithm stages. The detectors run every frame and do not take temporal information into account. This is instead evaluated at a higher level by the tracker, which matches new detections to existing tracks. Computationally expensive implementations of car, pedestrian and motion detectors and

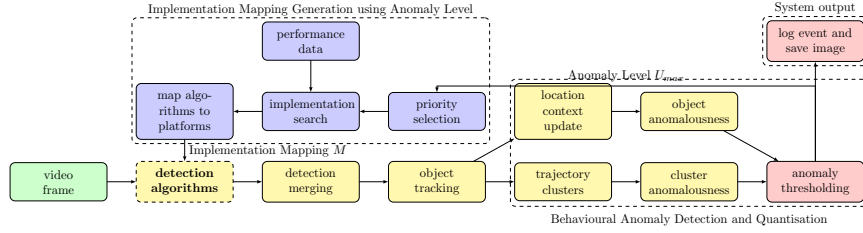


Figure 2: Flow diagram of frame processing, anomaly detection and algorithm mapping loop. Computationally expensive operations are in bold type. Boxes or groups with dotted lines show our contributions: #1(switchable hardware platforms for detection algorithms), #2(behavioural anomaly detection), #3(implementation mapping).

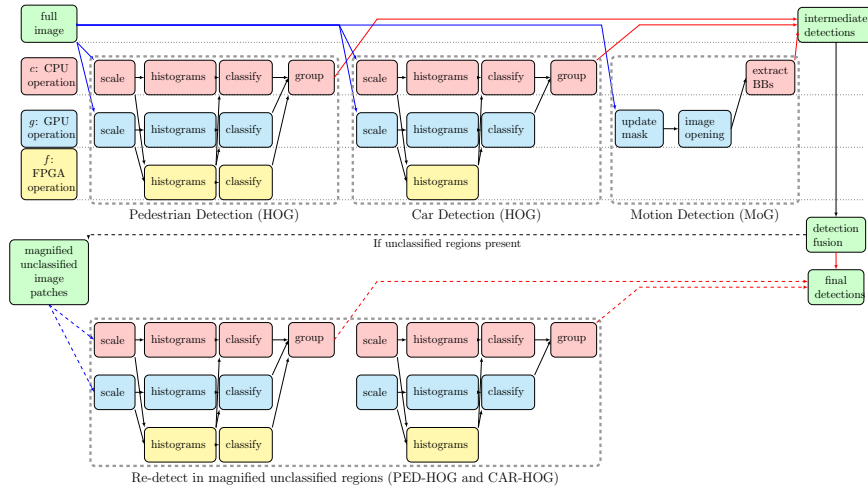


Figure 3: Hardware implementations of computationally expensive detection algorithms showing all possible mappings. Mnemonics for implementation stages are referred to in the text; HOG with image resizing on GPU, histogram extraction on FPGA and classification on GPU is labelled as *gfg*. The bottom half shows attempted reclassification of regions which are detected as only containing motion. This passes data through the same detectors as above.

Table 2: Total Resource Utilisation on Xilinx XC6VLX240 FPGA. This includes PCIe and DMA Logic and both Detectors: the Histogram extraction step for CAR-HOG and the Histogram and Classification steps for PED-HOG.

Resource	Resource Use	Percentage
Slice Registers	81888 of 301440	27%
Slice LUTs	77623 of 150720	51%
BlockRAM	217 of 832	26%
DSP48 Multiplier	108 of 768	14%

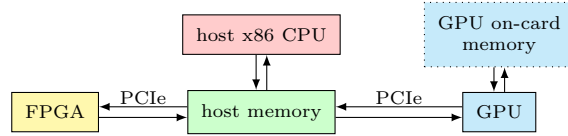


Figure 4: Desktop PC system with heterogeneous processors. Each processor is connected over PCI express and can access the host’s main memory. The FPGA and GPU each have private on-card memory.

algorithms run here. An expanded view of this stage is given in Fig. 3. These algorithms each had at least one accelerated version available; these are listed in Table 1.

### 3.1 Hardware Platform

The platform used contained a host CPU (2.4GHz dual-core Intel Xeon), an FPGA (Xilinx ML605 board with a XC6VLX240 device), and GPU (NVidia GeForce 560Ti). Data could be transferred via DMA (Direct Memory Access) between each processor over PCI express, as shown in Fig. 4. In the FPGA, the DMA and PCIe transfer logic ran at  $250MHz$  and the image processing application was clocked at  $160MHz$ . Resource use is given in Table 2. This table shows the total FPGA resources used to implement the PED-HOG and CAR-HOG detectors and the interface logic.

Application logic for each algorithm was autogenerated from Xilinx System Generator. There was sufficient capacity within the FPGA to hold the implementations in Fig. 3; because of this, partial dynamic reconfiguration was not required. Further platform details are given in [17].

### 3.2 Pedestrian Detection (PED-HOG)

Accelerated versions of the HOG pedestrian detector described by Dalal are used [11]. These are split into 3 tasks: (1) image resizing or scaling, (2) generation of gradient histograms in local cells and (3) block histogram generation, normalisation and SVM classification. These correspond to stages in the orig-



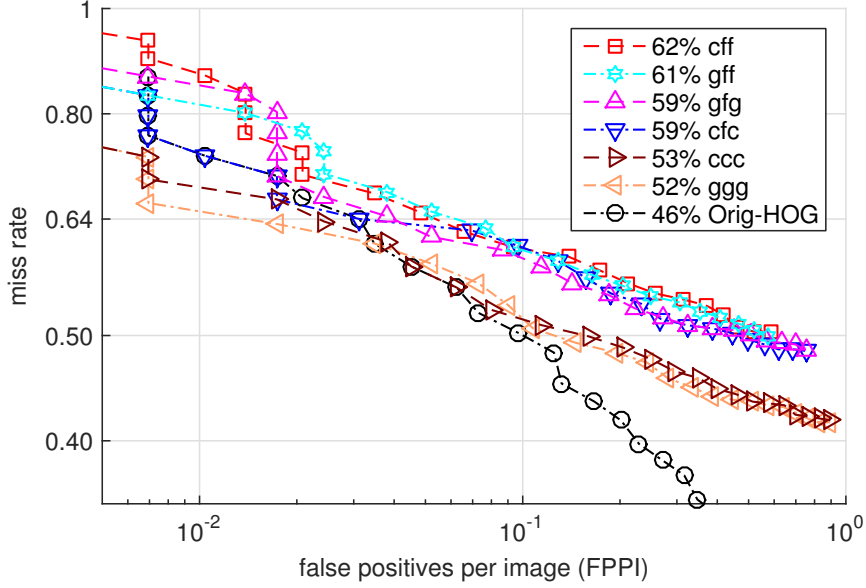


Figure 5: Detection error tradeoff curve for PED-HOG detector, showing false positives per image against miss rate.

inal algorithm [11]. A note on mnemonics: *xyz* refers, in general, to scaling on platform *x*, histograms on *y* and SVM classification on *z*. For example, an implementation which rescales the image on the GPU, generates histograms on the FPGA and classifies on the GPU is given the *gfg* mnemonic. Similarly, *ccc* means an implementation where scaling, histograms and classification is done on the host CPU. This is illustrated in Fig. 3, where the arrows on the left show the paths which can be taken through the algorithm. As Table 3 shows, the different implementations could be switched dynamically between frames with no loss of performance. Each version was trained on the INRIA pedestrians dataset [11]. Measurements of power, processing time and accuracy for each implementation when tested on the test portion of this dataset are given in Table 4. Detection accuracy for each version is shown in Fig. 5, and is comparable to the original algorithm.

While existing accelerated implementations of HOG are available, they do not take partitioning into account (i.e. do all the work on GPU), or use static partitioning (e.g. [16]). In addition to writing new implementations where nec-

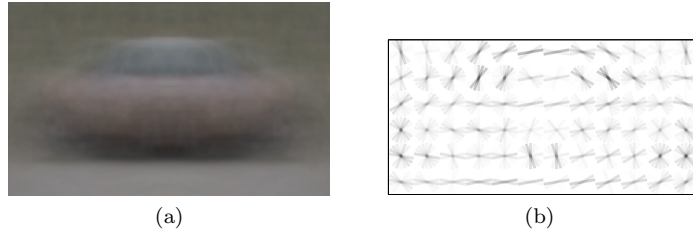


Figure 6: Training the CAR-HOG algorithm. (a) Composite image of all training samples. (b) The final oriented-gradients SVM car model.

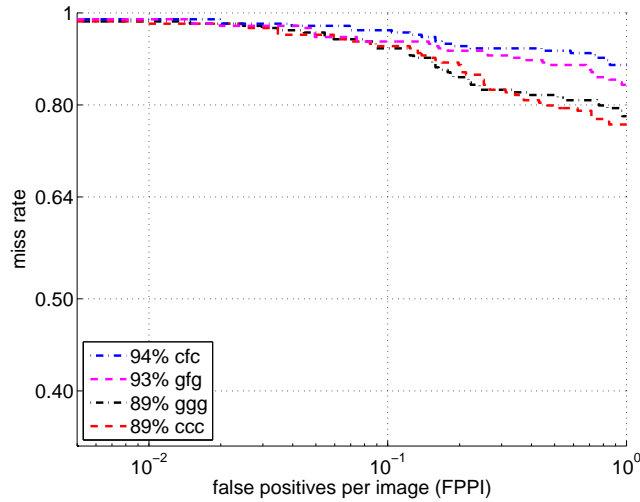


Figure 7: Detection error tradeoff curve showing false positives per image for CAR-HOG.

essary, our modification efforts for existing platforms (CPU and GPU) were focused on profiling (to identify the most appropriate stages to partition an algorithm and transfer intermediate data between processors) and interoperability; i.e. ensuring that cell histograms generated on FPGA could be transferred and classified on GPU, and that this produced acceptable results.

### 3.3 Car Detection (CAR-HOG)

Given that we had existing realtime implementations of HOG running on multiple architectures, these were modified to detect cars. The CPU and GPU versions of the HOG OpenCV code were modified using the parameters given by Dalal for vehicle detection [26]. Similarly, the *cfc* and *gfg* implementations

from [17] were modified to detect cars. In all cases, detection was performed at multiple scales in real time. While this results in a reduction in accuracy when compared to more sophisticated detectors, the implementation time required to produce multiple implementations of a different algorithm would have been prohibitive.

The car detector was trained using images from the 2012 PASCAL visual object classes challenge [27]. Image windows containing cars were extracted and resized so each car was 48 pixels high. Cars were selected if they were not marked as “occluded”, “truncated” or “difficult” in the ground truth and were originally  $\geq 40$  pixels high. 212 images were used, made up of original and horizontally flipped copies of 106 cars. The SVM model was trained as described in [11]. Fig. 6a shows a composite model of all positive training samples, and Fig. 6b shows the learned oriented-gradients model. Strong gradients on the vertical pillars are visible. Fig. 7 shows detector accuracy as an FPPI curve. This is less accurate than the pedestrian version; we speculate that this is due to the reduced set of training images and the wide variations in appearance of cars being viewed head-on and side-on. During testing, vans and lorries were also capable of being detected.

### 3.4 Motion Detection (MoG)

The Mixture of Gaussians (MoG) algorithm was used to perform background subtraction, thus segmenting foreground objects [28]. There were several motivations behind this; this technique allowed detection and then classification of small or distant objects, below the minimum size of the HOG detectors, and also detection of moving objects close together. The alternative to this would have been running both HOGs on a magnified version of the entire image, which would have not been possible with realtime processing. The OpenCV GPU implementation was used [29].

This was then morphologically opened (as both these steps were computationally expensive) before transferring to the host, where contour detection



Figure 8: Mixture-of-Gaussians motion detection algorithm on GPU. (a) Motion bounding boxes successfully identified. (b) False positives due to illumination changes and camera shake.

was performed and bounding boxes were produced. Every foreground region which was not yet classified was passed to additional pedestrian and vehicle detectors (as shown in the lower half of Fig. 3), so early identification of overlaps led to reduced processing at a later stage. An overlap criterion was used to compare bounding boxes; for pairs with  $\geq 90\%$  intersection, the smaller one was removed. Bounding boxes are defined as  $B = \{x_1 \dots x_2, y_1 \dots y_2\}$  and  $\text{area}(B) = (x_2 - x_1) \times (y_2 - y_1)$ . Thus  $B_i$  and  $B_j$  are compared and  $B_i$  is discarded if:

$$\frac{B_i \cap B_j}{\text{area}(B_j)} \geq 0.9 \ \& \ \text{area}(B_i) < \text{area}(B_j). \quad (1)$$

This performed well, as shown in Fig. 8a. Automatic gain correction or shake from the camera would occasionally incorrectly detect motion, as shown in Fig. 8b. In this case, all motion detection for that frame was discarded.

### 3.5 Detection Fusion & Object Tracking

The system had two sources of object detections: classified bounding boxes generated directly by CAR-HOG or PED-HOG, or unclassified regions of interest generated by the motion detector. Regions from the latter algorithm were extracted and magnified, then passed to each HOG algorithm again. This allowed most objects to be classified as human or vehicle, rather than simply an ob-

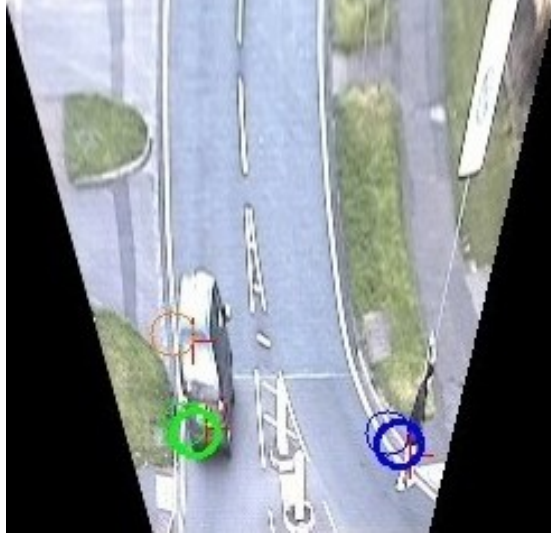


Figure 9: A transformed base plane image showing car (green circle) and pedestrian (blue circle) object trackers. Image is the lower part of the i-LIDS scene.

ject in motion. As we expect humans and vehicles to move at different speeds and frequent different regions in the scene, this allowed more accurate contextual scene information to be gathered than if we had simply detected moving objects without attempting to classify them.

To normalise inter-object distances and speeds, all detections were projected via an affine transform onto a base plane then smoothed with a constant-velocity Kalman filter, as Fig. 9 shows. This also matched per-frame detections to persistent trackers for each object; an elliptical distance measure was used to select a tracker to match a detection to, where the long axis of the ellipse pointed in the direction of travel of the tracker. When a tracked object was stationary, this became the Euclidean distance measure. This approach allowed unclassified detections to be matched to existing classified tracks. From this algorithm stage onwards, the volume of data to be processed was low enough that acceleration was no longer required, and all further processing was done on the host CPU.

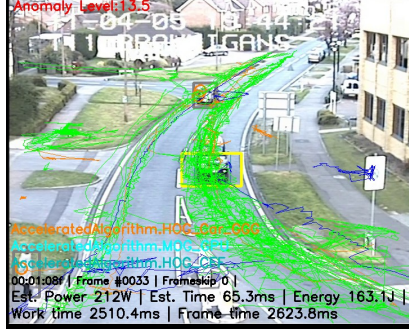


Figure 10: Object tracks grouped into trajectory clusters and re-transformed onto camera plane. Green, blue and orange tracks show cars, pedestrians and undetermined (motion-only) objects respectively.

### 3.6 Clustering Trajectories

To cluster tracked points into trees of trajectories, we re-implement a clustering algorithm described by Piciarelli and Foresti, who used it for detecting anomalies in traffic flow [7]. The original work looked at a motorway junction where fast-moving vehicles can take a small number of routes through the scene, but the scenes we apply this to are less structured. We consider urban scenes with two types of traffic participant, multiple entrances and exits, and areas where objects stop moving for long periods and may be considered part of the background. Our objective is to assign a trajectory  $T$  to one of a set of clusters  $\mathbf{C} = \{C_1, C_2, \dots, C_n\}$ . Each trajectory  $T_i = (t_0, t_1, \dots, t_n)$  represents tracked detections over several frames, where each  $t_i = (x_i, y_i)$ . Each cluster  $C_k = (c_0, c_1, \dots, c_n)$  is a vector of points  $c_j = (x_j, y_j, \sigma_j^2)$  with location and variance. Clusters are arranged in a tree structure and have zero or more children. A tree (starting with a root cluster) describes a single point of entry to a scene from one of the edges, and all observed paths taken through it from that point. Given a new (unmatched) trajectory  $T_u$ , all root clusters and their children are searched to a given depth.  $T_u$  is assigned to the closest  $C_k$  if the distance  $D$  is

below a threshold. This is done via:

$$D(T_u, C_i) = \min \left[ \frac{d(t_i, c_j)}{\sqrt{\sigma_j^2}} \right], j \in \{ \lfloor (1 - \delta)i \rfloor \dots \lceil (1 + \delta)i \rceil \} \quad (2)$$

where  $d(t_i, c_j)$  is the Euclidean distance between  $t_i$  in  $T_u$  and  $c_j$  in  $C_k$ , and  $j$  defines the range in  $C_k$  to search over. The lower and upper search bounds of the cluster points in  $C_k$  are governed by  $\delta = 0.4$ . Thus when matching long trajectories to longer clusters, more possible matches are allowed. This takes account of subsequent objects in one cluster not moving in exactly the same manner. Once a point is matched to a cluster, the corresponding cluster element  $c_j$  is updated by  $t_i$  with a learning factor  $\alpha = 0.05$ .

If no matching clusters are found, a new root cluster is created. As trajectories matched to a cluster are updated with new points, the points can be updated, the cluster extended, or new child clusters split off as points begin to diverge. Fig. 10 shows this approach on real data; several clusters representing pedestrian and vehicle motion can be seen.

### 3.7 Contextual Location Data

Contextual knowledge makes use of known information about the normal or common actions of participants within a scene. Simple features such as position and motion data can capture various anomalous behaviours including vehicles parked in an unusual location or those moving the wrong way down a road. Here we use unsupervised learning to learn scene context. We use the classified bounding boxes from the object and motion detectors during training sequences to build up two-dimensional per-pixel base plane heatmaps of object position or presence; these are shown in Fig. 11.

In a similar manner, Fig. 12 shows motion maps for pedestrians and vehicles in  $x$  and  $y$ . Velocity data was obtained from object trackers observed during the learning clips. For an object moving at  $\mathbf{v} = (v_x, v_y)$  and occupying  $(x_1, \dots, x_2, y_1, \dots, y_2)$ , the map  $\bar{\mathbf{v}}$  representing average per-pixel velocity is

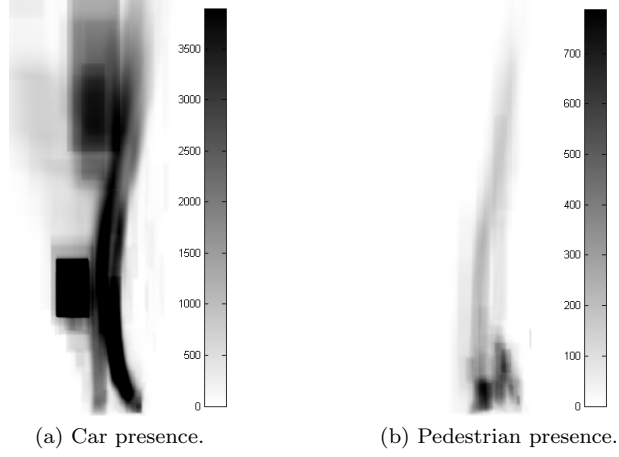


Figure 11: Presence intensity maps for cars and pedestrians in the PV3 dataset. The colour bar shows the average occupation rate of that pixel, and the  $x$  and  $y$  axes relate to the base-plane projection of the scene (partly shown in Fig. 9).

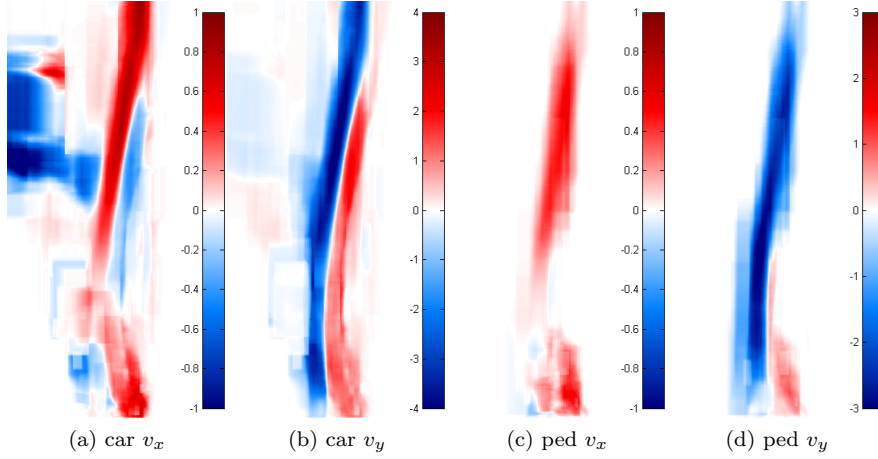


Figure 12: Base-plane motion maps built using learned movement from different object classes in PV3. In (b), on-road vertical motion away from (blue) and toward the camera (red) is distinct and clearly defined. The colour bar denotes average velocity at that pixel, in pixels per frame. The  $x$  and  $y$  axes correspond to the base plane image of the scene from Fig. 9.



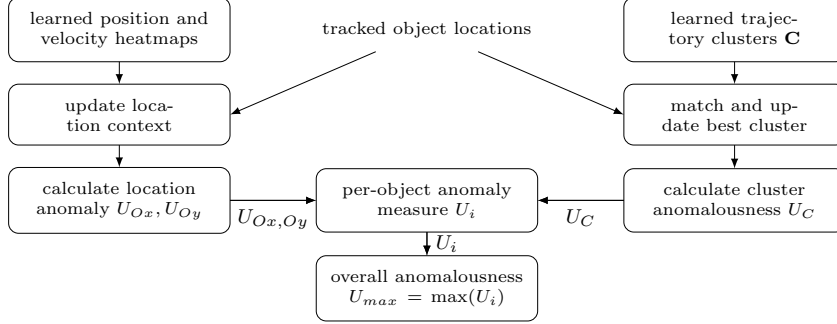


Figure 13: Flow diagram of contribution #2: anomaly detection process. Using tracked object locations identified by the detection algorithms, learned cluster, position and velocity information is used to assign anomaly levels to individual objects. Finally, an overall per-frame behaviour anomaly level is quantified.

updated:

$$\bar{\mathbf{v}}_{(x_1, \dots, x_2, y_1, \dots, y_2)} = (1 - \alpha)\bar{\mathbf{v}}_{(x_1, \dots, x_2, y_1, \dots, y_2)} + \alpha \mathbf{v} \quad (3)$$

Where  $\alpha = 0.0002$ . Fig. 12 shows motion maps for pedestrians and vehicles in  $x$  and  $y$ . This approach works well for most conceivable traffic actions. However, this does not capture more complex interactions between traffic participants. As this is unsupervised learning, errors from the object detectors will propagate to the heatmaps.

### 3.8 Detecting Anomalies

Using the algorithm in §3.6 and §3.7 we can define an anomalous object as one which is present in an unexpected area, or one present in an expected area but which moves in an unexpected manner. It follows that these objects will relate to events which are not present or under-represented in the training data and not representative of normal traffic flow. We use a Bayesian approach to determine if an object’s velocity in the  $x$  and  $y$  directions should be marked as anomalous, based on the average velocity  $\bar{v}$ .

First, we define an object anomaly measure  $U_O$  as the probability of a pixel

being associated with an anomaly  $p(A|D)$ , given a detection at that pixel:

$$U_O = p(A|D) = \frac{p(D|A)p(A)}{p(D|A)p(A) + p(D|\bar{A})p(\bar{A})}, \quad (4)$$

Where  $p(A)$  and  $p(D|A)$  are constants (given that the probability of detecting an anomaly at any point within the image is a constant, and the likelihood of detecting an object at any anomalous pixel is also constant). We set  $p(\bar{A}) = 1 - p(A)$  and we can, however, vary  $p(D|\bar{A})$ , which is based on the similarity between the observed motion  $\mathbf{v}$  and the learned mean per-pixel motion  $\bar{\mathbf{v}}$ , and is in the range (0.01,0.99). First we find  $d_v$ , a distance measure between  $\bar{\mathbf{v}}$  and  $\mathbf{v}$ . This is done separately for  $v_x$  and  $v_y$ .

$$\begin{aligned} \text{if } \quad & \text{sign}(\bar{v}) == \text{sign}(v) \ \& \ |v| > |\bar{v}|, \\ d_v \quad &= \text{sign}(\bar{v}) \times \max(W \times |\bar{v}|, |\bar{v}| + W) \\ \text{otherwise, } \quad & d_v = \text{sign}(\bar{v}) \times \min(-0.5 \times |\bar{v}|, |\bar{v}| - w). \end{aligned}$$

Where  $W$  and  $w$  are forward and reverse-directional constants. We perform linear regression, obtaining a line gradient of  $a = (0.01 - 0.99)/(d_v - \bar{v})$ . We then obtain an intermediate value  $k$  which is substituted into  $p(D|\bar{A})$ :

$$k = av + b, \quad (5)$$

$$p(D|\bar{A}) = \max(0.01, \min(0.99, k)). \quad (6)$$

$b$  is calculated in a similar manner to  $a$ . Finally, (4) is used to obtain  $U_{Ox}$ . This is repeated for  $U_{Oy}$ .

We combine this with  $U_C$ , an anomaly measure describing information about the abnormality of the current cluster associated with an object.  $U_C$  is based on one of two measures: transits and transitions. When an object moves from one cluster to any of its children or leaves the field of view, the number of *transits*

through that cluster is incremented. For a parent cluster  $C_p$  with children  $C_{c1}$  and  $C_{c2}$ , the number of trajectory *transitions* between the parent node and each of its children is also recorded. This builds up a frequency distribution between all children of any cluster. Anomalous trajectories can thus be identified. These metrics are combined; if  $C_i$  is a root node,

$$U_C(C_i) = \frac{1}{1 + \text{transits}(C_i)}, \quad (7)$$

otherwise if  $C_i$  is one of  $n$  child nodes of  $C_p$ ,

$$U_C(C_i) = 1 - \frac{\text{transitions}(C_p \rightarrow C_i)}{\sum_{j=1}^n \text{transitions}(C_p \rightarrow C_j)}. \quad (8)$$

The anomaly measure  $U_i$  for object  $i$  with an age of  $\tau$  frames can then be given:

$$U_i = w_o \frac{\sum_1^{\tau_i} U_{Ox}}{\tau_i} + w_o \frac{\sum_1^{\tau_i} U_{Oy}}{\tau_i} + w_c U_C(C_i), \quad (9)$$

Weights  $w_o, w_c$  are set such that two out of three anomaly indicators are needed to flag an object as anomalous.  $U_{Ox}$  and  $U_{Oy}$  are running totals averaged over  $\tau$ . For every frame, the overall anomaly measure is then  $U_{max} = \max(U_i)$ . This process is shown in Fig. 13. Taken together, objects which stop in areas where the normal behaviour is to move at speed in a particular direction are detected both by the object anomaly and the cluster trajectory algorithms (as a stationary object will still advance in time, out of the frequently-transited parent cluster and ultimately into a child cluster of its own). If an object is flagged as anomalous for a minimum time threshold  $t_A$ , its location and classification is logged and a snapshot of the video frame is saved.

### 3.9 Dynamic Algorithm Mapping

By this stage, we have reduced a video frame to  $U_{max}$ , a single scalar describing the anomaly level in that frame. We now use this to choose a mapping or set

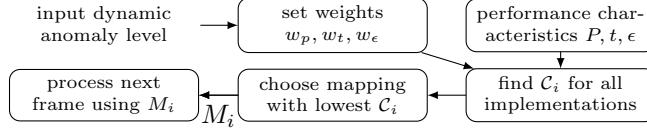


Figure 14: Flow diagram of contribution #3 - hardware mapping process. Using a dynamic anomaly level as input, the lowest-cost mapping is found using stored performance characteristics and used to process the next frame with.

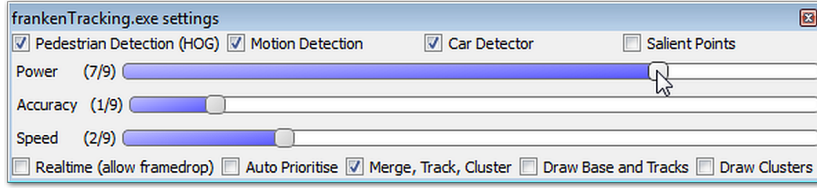


Figure 15: User- or algorithm-driven priority selection. Moving a slider to the right represents increased weight given to that performance characteristic. Other sliders are automatically moved to the left to compensate. See supplementary material for videos of this process.

of algorithm implementations  $M$  to process the next frame with. This process is shown in Fig. 14.  $M$  is recalculated every time a frame is processed, so new implementations can be selected in response to changing scene characteristics. Any mapping  $M_i$  must include one implementation of PED-HOG, CAR-HOG and MoG ( $m_{ped}, m_{car}, m_{mog}$ ) and can be any combination of the paths through Fig. 3. If a frame takes longer to process than realtime, subsequent frames are skipped to maintain realtime performance. Time spent processing every part of the algorithm in Fig. 2 is counted and used when calculating the number of frames which should be skipped. This includes time spent attempting to re-classify small unclassified regions as shown in the lower half of Fig. 2.

We assume that a higher level of anomaly in the scene should be responded to with increased processing resources to allow inference in a more timely fashion, and at the expense of power consumption. Time periods with low or zero anomaly levels will cause power to be reduced, resulting in longer processing times, lower accuracy and more frames being dropped. We define  $R = \{w_p, w_t, w_\epsilon\}$  as the current prioritisation setting, with each  $w$  being individual prioritisations for time, power, and miss rate. Ten credits are allocated

between the three weights; the use of integer weights was a limitation of the user interface library.  $R$  is set by assigning a given fraction of credits to each weight. This can be done manually (by always assigning the maximum weight possible to speed or power; see Fig. 15 for an example) or automatically, by maximising speed when  $U_{max}$  is above a threshold or power when below. (Some hysteresis is built in by using two separate thresholds.)

Using the calculated priorities and set of desired algorithms, implementation mapping is then run.  $M$  is selected by choosing the lowest-cost implementation. For a given mapping  $M_i$ , a cost function is used:

$$C_i = w_p P_i + w_t t_i + w_\epsilon \epsilon_i, \quad (10)$$

Here, the implementation performance characteristics  $P, t$  and  $\epsilon$  represent system power consumption, frame processing time, and detection accuracy expressed as miss rate. These costs are incurred while processing a frame; that is, running  $m_{ped}, m_{car}, m_{mog}$  in  $M_i$ . These are calculated from the values shown in Table 4.  $t_i$  is the sum of the processing time of  $M_i$ .  $P_i$  is the average power consumed while processing  $m_{ped}$ ,  $m_{car}$  and  $m_{mog}$  (i.e.  $(P_{car}t_{car} + P_{ped}t_{ped} + P_{mog}t_{mog})/t_i$ ).  $\epsilon_i$  is a ranked measure of the miss rate of  $m_{car}$  and  $m_{ped}$ . A  $M$  with *ped-ggg*, *car-ggg*, *mog-gpu* would have  $t_i = 60\text{ms}$ ,  $P_i = 225W$  and  $\epsilon_i$  as a ranked accuracy measure. The points in Fig. 16 show discrete  $t_i$  and  $P_i$  calculations for every possible mapping.

Normalisation of the  $P, t, \epsilon$  values in (10) is not required. The calculated values can be used directly, as any decision on normalisation is ultimately application specific and represents the relative priority designers or users would assign to longevity of operation of any system vs. fast detection capability. The normalisation approach used (including use of integer weights for  $R$ ) balances all of these factors and ensures that multiple mappings are used where appropriate.

This tradeoff is shown graphically in Fig. 17, where points in stronger red, green or blue signify the majority of algorithm stages being mapped to FPGA,

GPU or CPU respectively. The lowest-cost mapping is then selected and used to process the next frame and generate new detections.

## 4 Methodology and Results

Having described the system and the parameters of interest (namely  $R$  and its effect on  $M$ ), we now test its ability to detect parked vehicles on two video datasets of scenarios. These are described here, followed by our methodology.

### 4.1 Scenarios

We use two representative scenarios of anomalous behaviour in an urban environment, both gathered from a static camera. The *BankSt* scenario involves a static camera monitoring a busy four-way crossroads in daylight, at  $720 \times 400$  and  $30fps$  with little camera shake present; see Fig. 1b. Vehicles are often stationary for extended periods while waiting at the traffic lights in both directions. BankSt has a 18-minute training and 4 minute test video. The training data contains no parked vehicles, only normal traffic patterns. The testing data contains one event where a vehicle parks in a forbidden location.

The i-LIDS dataset [30] is a UK government video dataset of various indoor and outdoor scenarios, used for evaluation of anomaly detection algorithms. We use a Parked Vehicle detection scenario, PV3, for real-time evaluation. This consists of 5.5 hours each of training and testing footage of a scene as shown in Fig. 18b. Ground truth of every event is provided as a timestamped description. The video is in colour, visible-light range and is  $720 \times 576$  at  $25fps$ . It overlooks an urban road; several turnings are present on the left and right and traffic often queues to enter a roundabout beneath the frame. Video clips from several months apart, in various weather conditions and at different times of day or night are provided. Strong shadows, camera shake, camera movement in between clips, videotape artefacts and noise interfere with reliable scene analysis. PV3 is poorer quality and considerably more challenging than our BankSt data. This

Table 3: Frame processing time  $t_{proc}$  for frame  $n$  using the HOG implementation listed, where the previous frame ( $n - 1$ ) was processed with the same (column 2) or different (column 3) implementation. Column 4 shows times for “different”  $t_{proc}$  as a percentage of the “same”  $t_{proc}$ . Image size is  $1024 \times 768$ . Implementation acronyms defined in §3.2.

Impl.	$t_{proc}$ (ms) for given impl. when previous impl. was:		Time for $a \rightarrow b$ , as % of $a \rightarrow a$
	$a \rightarrow a$ (same)	$a \rightarrow b$ (different)	
ccc	776.6	765.1	98.5
ggg	55.5	55.7	100.3
gfg	119.1	123.1	103.3
cfc	646.4	640.1	99.0
cff	96.6	95.1	98.4
gff	85.5	88.2	103.1

is apparent in the comparison shown in Fig. 18. It is, however, instructive as an example of the data that anomaly detection algorithms may have to work with in real-world tasks.

## 4.2 Methodology

The processing platform described in previous work [17] was first evaluated to ensure that no substantial delay was present when differing implementations were selected. Frame processing times were recorded when either maintaining constant mappings or switching them between subsequent frames. Switching between different mappings involved changes on up to three platforms. FPGA implementation switching involved setting or clearing bits to select a different processing pipeline; this was done as part of the command to start a DMA transfer and involved no reconfiguration or overhead. CPU and GPU implementation switching involved taking different software branches while a frame was being processed. This is included in the processing time, so any switching costs were thus negligible. As Table 3 shows, the time difference when switching mappings was always under 4% of the non-switching time. This ensured that (10) did not need to include a switching cost.

In each test we performed, we compare the performance of three prioritisation methods with the following labels:

**speed** speed manually prioritised, auto prioritisation off;

**power** power manually prioritised, auto prioritisation off;

**auto** anomaly-controlled auto prioritisation on.

For both datasets we manually register points in the camera image onto a base plane for a homography transform. Unlike in [9], this is the only manual intervention required; we do not restrict the detection of anomalies to a known area in the image through masking, but allow them to occur anywhere. This also allows expansion to multiple scenarios.

The clusters and heatmaps are then learned by running a training clip containing no anomalous behaviour, 17-20 minutes in length. The same clusters and heatmaps were used for all tests. Following Loy’s approach [4], we define an anomalous event as “observed behaviour which is absent or rarely present in the training data”. Thus vehicles parking in forbidden areas are used as anomalous events; we define these as objects which have  $U_i$  over a set threshold for a time period  $t_A$ .

We first test our system on BankSt to demonstrate extensibility to a variety of surveillance scenarios, then perform a complete evaluation on i-LIDS. The single event in the four minute BankSt test clip was detected using all three prioritisation modes (as shown in Fig. 1b), along with one false positive, in each case caused by stationary vehicles at the traffic lights. More analysis is performed on i-LIDS as it is a larger and more challenging scenario. It is also an open dataset which has benchmarks and results already reported in the literature. The default i-LIDS criteria require that an anomaly event must persist for 60 seconds before it can be logged, then it must be recorded within 10 seconds following this to count as a true positive. As our algorithm only needs 10–15 seconds to detect events, events are treated as true positives if the recorded start time is within 70 or 75 seconds after the ground-truth time



Table 4: Performance details for algorithm implementations on  $770 \times 578$  video. Processing Time (ms), Overall System Power Consumption (W) and Detection Accuracy (log-average miss rate (%)) shown. Baseline power consumption was 147W.

Algorithm	Platforms	Time (ms)	Power (W)	Miss rate (%)
PED-HOG	<i>ggg</i>	17.6	229	52
	<i>cff</i>	23.0	190	62
	<i>gff</i>	27.5	186	61
	<i>gfg</i>	39.0	200	59
	<i>cfc</i>	117.3	187	59
	<i>ccc</i>	282.0	191	53
CAR-HOG	<i>ggg</i>	34.3	229	89
	<i>cfc</i>	175.6	189	94
	<i>gfg</i>	60.0	200	92
	<i>ccc</i>	318.0	194	89
MoG	GPU	8.1	202	N/A

recorded when the vehicle actually parks. However, we require anomalous events to be localised to the object which is parked wrongly; this is more discriminative than the original i-LIDS criteria, which only require a binary alarm signal in the presence of an anomaly. We evaluate with the anomaly detection window  $t_A$  set to 10 and 15 seconds, log detected events, and compare them with ground-truth event data.

### 4.3 Results

We first summarise object detection performance by considering each measurement (power, speed, accuracy) separately. This allows us to gather data on the individual implementation performance characteristics required for generating  $M$ . We then discuss anomaly detection on i-LIDS.

#### 4.3.1 Object detection performance characteristics

Performance for object detection is given in Table 4. Here, the pedestrian detectors are always more accurate than the car versions. The GPU-based detectors are faster and consume more power, while implementations which perform more

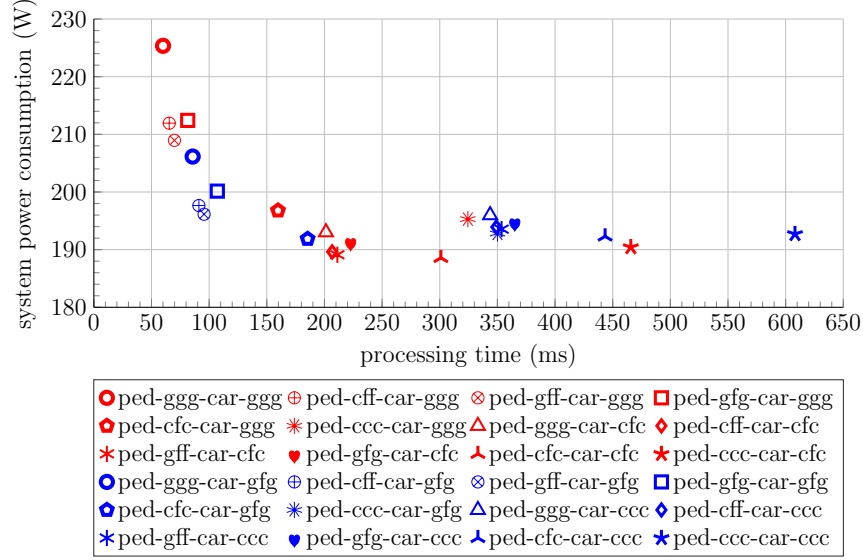


Figure 16: Design space exploration of power consumption vs. processing time for every car, pedestrian and motion detector implementation. All points include motion detection on GPU.

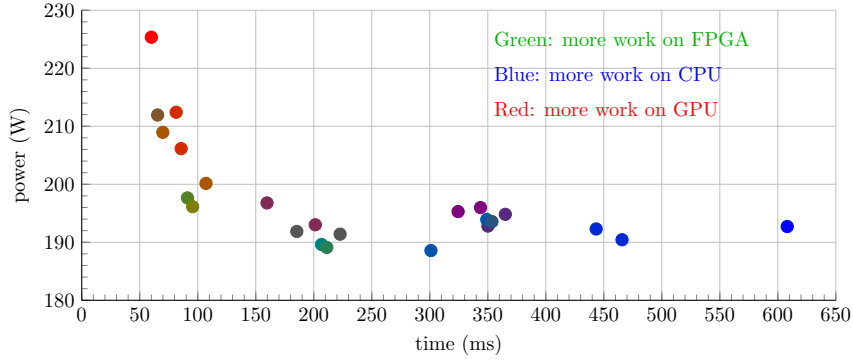


Figure 17: Power and time plots of all possible solutions, where each solution consists of one CAR-HOG, PED-HOG and MoG detector, as described in 4. A greener dot represents a solution with most operations mapped to FPGA, while bluer and redder dots represent those with most operations mapped to CPU or GPU respectively.

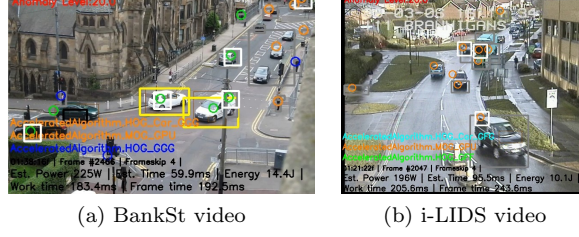


Figure 18: Examples of video quality and effect on detection ability: in BankSt (a), classification of most objects as either pedestrian (blue circle) or car (green circle) is possible. In i-LIDS, (b), many tracks remain as unclassified moving objects (orange circles).

processing on FPGA have reduced power consumption and accuracy.

There are several causes of the accuracy differences seen in the final column of Table 4; these are due to algorithm and platform differences between processors. Our FPGA implementation is the most inaccurate, due to numerical and algorithm simplifications. Fixed-point arithmetic must be used on FPGA, as this results in significantly lower resource use. This causes errors in precision to accumulate when compared to single-precision floating point versions seen on CPU and GPU. The FPGA histogram generation step omits the Gaussian weighting of pixel blocks at a slight cost in accuracy. The FPGA classification stage also uses a simplified block normalisation step compared to the software implementations (L1-normalisation followed by a square-root, compared to two-stage L2 normalisation present in the original algorithm; see [26]). Both of these decisions were taken to reduce the FPGA resource use and minimise expensive calculations such as division operations. Discrepancies between the CPU and GPU accuracy measurements are caused by execution of floating-point operations in a different order to achieve maximum parallelisation, and slight differences in the pixel cell generation code. These per-platform differences, when combined, account for the variations seen in miss rate. More details are given in [17].

Power and time data was obtained while playing a  $770 \times 578$  video and the detection accuracy figures refer to the testing portion of the dataset used

for training (either INRIA or PASCAL). Power measurements in this table were gathered using a plug-in meter which recorded total system power consumption. We plot this power-time tradeoff in Fig. 16 (showing exact placement of each mapping solution) and in Fig. 17 (showing the proportion of each work carried out on a particular processor). Here the solutions form a Pareto curve, stretching from  $\{ped-ggg, car-ggg, mog-gpu\}$  at top left via  $\{ped-gff, car-gfg, mog-gpu\}$  to  $\{ped-cfc, car-cfc, mog-gpu\}$  in the lower centre, with other points shadowed. If a single mapping was both lowest-power and fastest, it would overshadow all other mappings and always be the optimal mapping to use in all situations; this is not the case.

Due to the design choices made when producing implementations for each platform, the ranking of accuracy measurements between implementations is fixed and does not depend on the choice of dataset used to train or test the individual detectors; there will be no dataset where e.g. *gff* is more accurate than *ggg*. The cost function expresses the tradeoff of choosing a less accurate implementation *of the same algorithm evaluated on the same dataset* over one which uses more power; thus, the performance information in Table 4 is not data-dependent. In support of this conclusion, Dollar *et al.* showed that ranking of object detector performance remains relatively consistent between different datasets [12].

#### 4.3.2 Anomaly detection performance characteristics

All detected anomalous events are logged and a snapshot taken. We compare this to the ground truth using both precision  $p = TP/(TP + FP)$  and recall  $r = TP/(TP + FN)$  measures, where  $TP$ ,  $FP$ ,  $FN$  are true positives, false positives and false negatives respectively. The  $F_1$ -score is also calculated:

$$F_1 = (\alpha + 1)rp/(r + \alpha p) \quad (11)$$

Table 5: Detection of parked vehicle events on i-LIDS PV3 for each prioritisation mode. True positives (TP), false positives and negatives (FP,FN), precision (p), recall (r) are shown.  $F_1$ -scores shown for operational awareness ( $OA$ ) and event logging ( $EL$ ).

Priority	TP	FP	FN	$p$	$r$	$F_{1,OA}$	$F_{1,EL}$
for $t_A = 10$ seconds							
power	4	44	26	0.083	0.133	0.0957	0.1317
speed	6	51	25	0.105	0.194	0.1254	0.1913
auto	6	53	25	0.102	0.194	0.1226	0.1912
for $t_A = 10$ seconds, daylight only							
power	4	29	23	0.121	0.148	0.1294	0.1475
speed	6	40	22	0.130	0.214	0.1510	0.2118
auto	6	42	22	0.125	0.214	0.1466	0.2115
for $t_A = 15$ seconds							
power	2	15	29	0.118	0.065	0.0910	0.0652
speed	8	13	23	0.381	0.258	0.3259	0.2594
auto	4	14	26	0.222	0.133	0.1797	0.1342
for $t_A = 15$ seconds, daylight only							
power	2	10	29	0.167	0.065	0.1067	0.0652
speed	8	8	23	0.500	0.258	0.3752	0.2601
auto	4	10	26	0.286	0.133	0.2030	0.1345

Table 6: Processing performance for power, speed and time prioritisations. Frame skip percentage( $fskip$ ), processing time with and without overheads ( $t_{all}$  and  $t_{work}$ ) compared to source frame time ( $t_{src}$ ), total estimated and mean estimated power above baseline ( $P_{work}^*$ ) are shown. Baseline (idle) power consumption was 147W.  $F_1$  accuracy scores are also shown.

Priority	$fskip$ (%)	$t_{all}/$ $t_{src}$ (%)	$t_{work}/$ $t_{src}$ (%)	$P_{work}^*$ (W)	$F_{1,OA}$
power	75.8	125.4	87.9	49.1	0.0910
speed	59.5	124.3	81.7	72.8	0.3259
auto	66.5	127.8	83.4	61.9	0.1797
daylight only					
power	75.5	125.4	87.9	49.1	0.1067
speed	60.1	124.2	82.0	78.4	0.3752
auto	66.0	122.0	83.4	61.8	0.2030

where  $\alpha$  is a “recall bias” measure, set to 0.55 for real time operational awareness (so that false alarm rate is reduced and operator confidence is maintained), or 60 for event logging (so that all plausible events are logged), giving  $F_{1,OA}$  and  $F_{1,EL}$  respectively. Details of these are given in Table 5 for both  $t_A$  measures, for the two manual and one automatic prioritisation modes. There are no ground-truth events in the night clips but these generate a large proportion of false positives, so we also show results for daylight-only clips (those labelled “day” or “dusk”). From this table, prioritising for speed always improves accuracy, but prioritising reduced power consumption reduces accuracy. Automatic prioritisation provides a compromise between these two extremes. Extending the  $t_A$  to 15 seconds allows a significant reduction in the number of false positives.

In Table 6 we show performance details for all prioritisation modes. The  $t_{all}$  column shows per-frame execution time including overheads as a percentage of the time available per video frame,  $t_{src} = 40ms$ . The processing time column  $t_{work}$  does not include overheads (i.e. assumes that a raw frame is already in memory and annotated per-frame output is not needed so that only events are logged). The system runs faster than real time here. When running individually, some implementation mappings are able to process every frame in real time, as

shown by the processing times in Table 4. Algorithms are run sequentially, so for 40ms frames we can run one or two of the detection algorithms for every frame, but not all three. The overall system is, however, able to process a video stream in real time by skipping a proportion of frames.

The percentage of skipped frames is in the *fskip* column; as expected, when we optimise for speed then frames are processed faster so fewer frames are dropped. From this table, power prioritisation reduces overall power consumption while speed increases it, showing that the prioritisation modes have some effect and behave as expected. The ‘power’ row in Table 6 maps most processing to FPGA, so here 75% of frames are skipped. When speed is prioritised, the ‘speed’ prioritisation maps everything to GPU, so when all work is done on GPU then only 59% of frames are skipped. For ‘auto’, the implementations used vary so an average frameskip value is shown. Fast algorithms which skip fewer frames have higher accuracy, as shown by the relationship between the *fskip* and  $F_1$  values in Table 6.

Example detections, logged after  $t_A$  seconds, are given in Fig. 19. While true positives are detected in a variety of conditions (including at dusk), even in the best case only up to 50% of events are detected. There are also many false positives and negatives, due mainly to shortcomings in the object or motion detectors. False positives can be caused by the MoG algorithm flagging patches of empty road as foreground (Fig. 19g), but mitigating this may cause slow-moving or stationary traffic in that region to be ignored. Errors can also be caused by limitations of the HOG implementations (such as in Fig. 20), or failure to detect occluded objects (Fig. 19h). Video quality has an impact here too; as Fig. 18 shows, higher-quality video allows more reliable classification (into pedestrian and car classes) of moving objects. Other limitations such as in Fig. 19i are also apparent; as we have failed to associate the anomalous event with the object causing it, this incident counts as *two* errors (FN and FP).

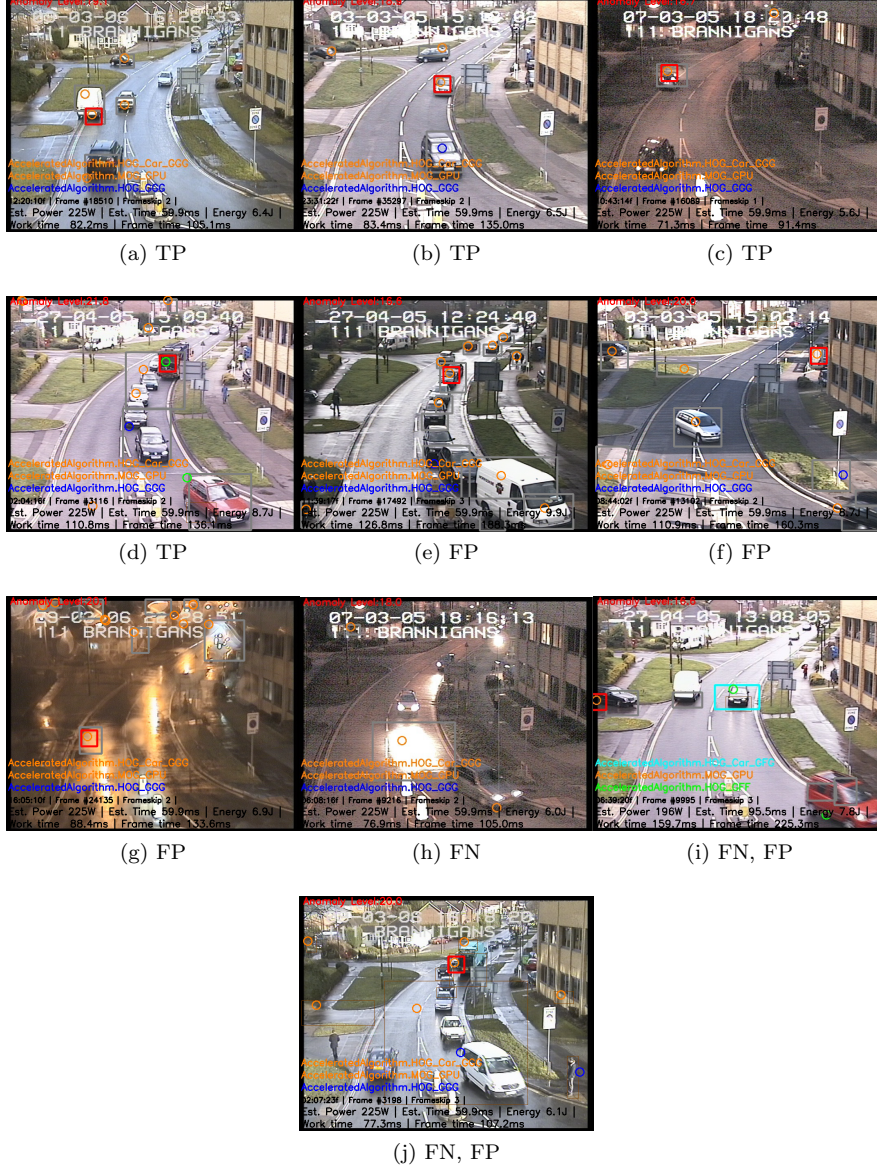


Figure 19: Examples of positive detections and failure modes of anomaly detector on i-LIDS. Anomalies are marked by red boxes. (a)–(d): true positives of various objects in different locations in quiet and busy surroundings. (e)–(g): false positives from slow-moving traffic, an object moving off-screen and incorrect background subtraction respectively. (h): a car is occluded behind the road sign on the right; this is treated as a false negative. (i) is a false negative and a false positive: the anomaly detector identifies the car on the left instead of the van parked beside it. (j) is an anomaly identified outside the allowed time window, so counts as a false negative and false positive.



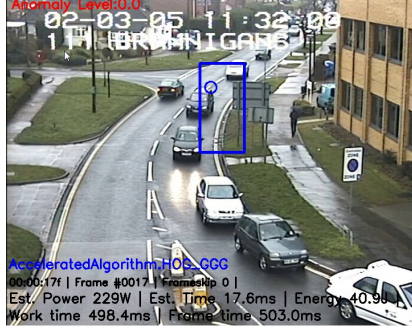


Figure 20: Failure mode of object detection stage, showing a false positive (blue rectangle representing a pedestrian). These affected the measured accuracy, both directly and through the cluster and object abnormality measures.

## 5 Analysis

Here we concentrate on the most significant results from these experiments and compare them to existing work. Considering the tradeoffs between power consumption, accuracy and time, we note that in this case we are constrained by the real-time processing requirement. Power and accuracy are therefore the main characteristics we can vary, as we have limited flexibility over the time requirement.

The key results are found in Tables 5 and 6. From Table 6, there is a  $29W$  range in average power consumption above baseline. The platform used was based on a desktop PC and so was not optimised for low power consumption (hence the relatively high  $147W$  idle power). With the system run with speed prioritised and the FPGA disconnected, power consumption  $P_{avg}$  was  $208W$ , or  $62W$  above baseline.  $P_{avg}$  consumed under auto-prioritisation with FPGA enabled was  $61.9W$ , but this is specific to this dataset; video datasets with fewer moving objects or longer idle periods would have a considerably lower  $P_{avg}$  than this. This is also because of the lack of *gff* and *cff* implementations for CAR-HOG, which would reduce  $P_{avg}$  further.

We consider  $P_{avg}$  vs. accuracy in Fig. 21. Using *auto* prioritisation allows a 10% increase in accuracy over the *power*-optimised option, at a cost of  $12W$  extra in  $P_{avg}$ . A further 17% improvement in  $F_1$ -score (moving from *auto* to

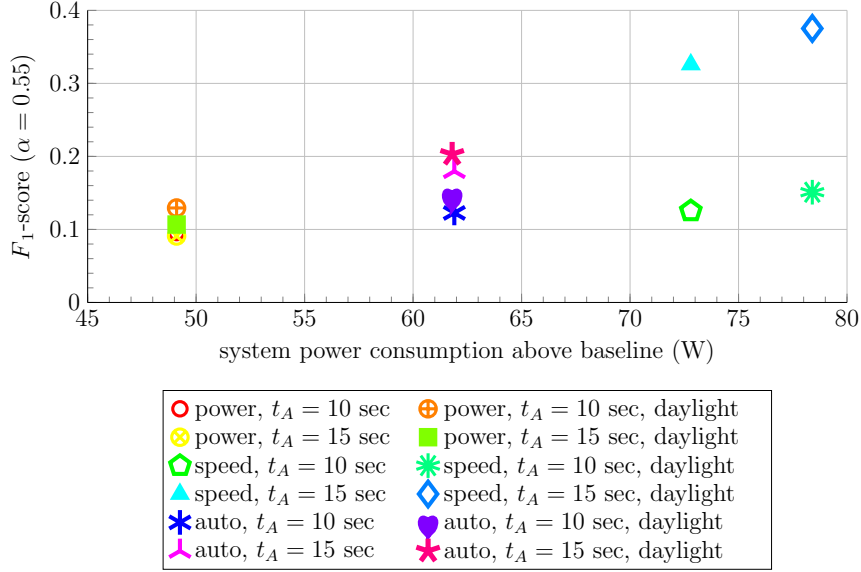


Figure 21: Error rate vs. power consumption:  $F_1$  operational awareness scores ( $\alpha = 0.55$ ) against power consumption for various time thresholds, for various lighting conditions.

*speed*) costs an extra 17W. When compared to *speed*, the 12W power reduction moving to *auto* reduces accuracy by 45% of baseline, while the fully-optimised *power* option loses 72% accuracy for 32% in power savings from best-case. This is most apparent for longer detection windows, on daylight-only clips with higher levels of anomalous behaviour. These measurements and Fig. 21 itself show a clear relationship between power consumption and overall detection accuracy.

## 5.1 A Comparison to State-of-the-Art Work

Following other researchers' definitions of anomalous events (events not present in and different to the training data) [4], we argue here that detection of parked vehicles is a reasonable representation of the more general task of detecting anomalous events in video. The only previous work to perform parked vehicle detection on the full i-LIDS PV3 dataset has been by Albiol *et al.* [9]. They used spatiotemporal maps and manually applied lane masks in each clip to signify which scene regions could have valid detections. Using this approach, they reached  $p$ - and  $r$ - values of 100% in some clips. Performance information is

not given, but video frames are downsampled to  $320 \times 240$  to decrease evaluation time. They note that, as an approach for detecting slow-moving and stationary vehicles, background subtraction has various limitations. They also encountered similar difficulties as those shown in Fig. 18. Similar work by Lee *et al.* performs illegal parked vehicle detection on short clips extracted from the i-LIDS set; they achieve close-to-realtime performance at original resolution, and detect all events in the clips with no false positives [10]. Thus, if we only take into account accuracy measurements, we are unable to improve on these existing results. Here, however, we consider a novel and different problem: that of automated power-aware anomaly detection rather than monitoring lane occlusion in a manually masked region. The only manual intervention we required in these experiments was to register points in each camera viewpoint to perform an affine transform onto the base plane. In future, this step could be done automatically.

The detection algorithms in this work are the most time- and power-intensive components of this system. Accuracy in these could be improved by implementing more sophisticated algorithms on one or more accelerated platforms. However, this would require considerable development time if each version were to be implemented. As Dollar *et al.* note, many of the current most accurate pedestrian detectors are based on HOG, so the hardware acceleration techniques documented in this work and the overall conclusions would apply if more accurate versions were substituted in the future [12]. As the Pareto curve in Fig. 17 shows, any such implementation with known accuracy would always involve a compromise between power and speed too; each measurement may become a priority at any point in time.

The other work concerning power-aware resource allocation in video is by Llamocca and Pattichis [24]. This paper is concerned with tasks at higher levels of abstraction than their work, but we are able to use the calculated anomaly level to dynamically drive selection of the optimal mapping, rather than relying on a user-selected accuracy level. This results in a fully automated system. We are thus able to tie low-level performance constraints such as power

consumption/energy to high-level ‘accuracy of detected behavioural events’.

## 6 Conclusion

We have described a real-time system which performs parked vehicle detection along with classification of humans and cars. This can select between various algorithm implementations on a mixture of FPGA, GPU and CPU, each of which have different power consumption, implementation runtime and algorithm accuracy characteristics. These are selected dynamically based on a feedback loop driven by the number of objects in a video perceived to be behaving anomalously (in this case, parking in forbidden areas). This allows us to dynamically trade off power consumption against detection accuracy and shows benefits when compared to fixed power-or speed-optimised versions. We evaluated this on a smaller dataset and performed a full characterisation on the larger i-LIDS PV3 dataset. This shows a clear link between processing power consumption and event detection accuracy; compared to power-prioritised selection, automatic anomaly-driven mapping is 10% more accurate but draws 12W more power.

[Supplementary Material] Supplementary material involving screen capture videos of the system performing detection on surveillance sequences, and showing priorities being adjusted as in Fig. 15, are available at <http://youtu.be/PSF-12a8Lo0> for BankSt and [http://youtu.be/Xm2bd0\\_TfnM](http://youtu.be/Xm2bd0_TfnM) for i-LIDS. Code and BankSt videos are available at <http://homepages.ed.ac.uk/cblair2/csvt/>.

## Acknowledgement

We gratefully acknowledge the work done by Scott Robson at Thales Optronics for his FPGA implementation of the histogram generator for the CAR-HOG detector. Thanks also to John Thompson for helpful comments on the manuscript.

This work was supported by Thales Optronics, the Engineering and Physical Sciences Research Council (EPSRC) Grant number EP/J015180/1 and the MOD University Defence Research Collaboration in Signal Processing.

## References

- [1] R. Benenson and M. Mathias, “Pedestrian detection at 100 frames per second,” in *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, 2012, pp. 2903–2910.
- [2] M. Hahnle, F. Saxen, M. Hisung, U. Brunsmann, and K. Doll, “FPGA-Based Real-Time Pedestrian Detection on High-Resolution Images,” in *Computer Vision and Pattern Recognition Workshops (CVPRW), 2013 IEEE Conference on*. IEEE, Jun. 2013, pp. 629–635.
- [3] D. Bacon, R. Rabbah, and S. Shukla, “FPGA Programming for the Masses,” *Queue*, vol. 11, no. 2, p. 40, Feb. 2013.
- [4] C. C. Loy, T. Xiang, and S. Gong, “Detecting and discriminating behavioural anomalies,” *Pattern Recognition*, vol. 44, no. 1, pp. 117–132, Jan. 2011.
- [5] B. T. Morris and M. M. Trivedi, “A Survey of Vision-Based Trajectory Learning and Analysis for Surveillance,” *IEEE Trans. Circuits and Systems for Video Technology*, vol. 18, no. 8, pp. 1114–1127, Aug. 2008.
- [6] S. Atev, O. Masoud, and N. Papanikolopoulos, “Learning Traffic Patterns at Intersections by Spectral Clustering of Motion Trajectories,” in *2006 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, Oct. 2006, pp. 4851–4856.
- [7] C. Piciarelli and G. Foresti, “On-line trajectory clustering for anomalous events detection,” *Pattern Recognition Letters*, vol. 27, no. 15, pp. 1835–1842, Nov. 2006.

- [8] B. T. Morris and M. M. Trivedi, “Learning, modeling, and classification of vehicle track patterns from live video,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 9, no. 3, pp. 425–437, 2008.
- [9] A. Albiol, L. Sanchis, A. Albiol, and J. M. Mossi, “Detection of Parked Vehicles Using Spatiotemporal Maps,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 12, no. 4, pp. 1277–1291, Dec. 2011.
- [10] J. Lee, M. S. Ryoo, M. Riley, and J. Aggarwal, “Real-Time Illegal Parking Detection in Outdoor Environments Using 1-D Transformation,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 19, no. 7, pp. 1014–1024, Jul. 2009.
- [11] N. Dalal and B. Triggs, “Histograms of Oriented Gradients for Human Detection,” in *Proceedings of the 2005 IEEE Conference on Computer Vision and Pattern Recognition (CVPR’05)*. IEEE, 2005.
- [12] P. Dollár, C. Wojek, B. Schiele, and P. Perona, “Pedestrian Detection: An Evaluation of the State of the Art,” *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 34, no. 4, pp. 743–762, Jul. 2011.
- [13] P. Dollár, R. Appel, S. Belongie, and P. Perona, “Fast feature pyramids for object detection,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 8, pp. 1532–1545, Aug. 2014.
- [14] V. Prisacariu and I. Reid, “fastHOG-a real-time GPU implementation of HOG,” Department of Engineering Science, Oxford University, Tech. Rep. 2310, 2009.
- [15] S. Martelli, D. Tosato, M. Cristani, and V. Murino, “FPGA-based pedestrian detection using array of covariance features,” in *Distributed Smart Cameras, ACM/IEEE International Conference on*, 2011, pp. 1–6.

- [16] S. Bauer and S. Kohler, “FPGA-GPU architecture for kernel SVM pedestrian detection,” in *Computer Vision and Pattern Recognition Workshops (CVPRW), IEEE Conference on*. IEEE, 2010, pp. 61–68.
- [17] C. Blair, N. M. Robertson, and D. Hume, “Characterising a Heterogeneous System for Person Detection in Video using Histograms of Oriented Gradients: Power vs. Speed vs. Accuracy,” *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, vol. 3, no. 2, pp. 236–247, 2013.
- [18] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, “Object detection with discriminatively trained part-based models,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 9, pp. 1627–45, Sep. 2010.
- [19] B. Cope, P. Y. K. Cheung, W. Luk, and L. Howes, “Performance Comparison of Graphics Processors to Reconfigurable Logic: A Case Study,” *IEEE Trans. Computers*, vol. 59, no. 4, pp. 433–448, Apr. 2010.
- [20] Q. Wu, Y. Ha, A. Kumar, S. Luo, A. Li, and S. Mohamed, “A heterogeneous platform with GPU and FPGA for power efficient high performance computing,” in *2014 International Symposium on Integrated Circuits (ISIC)*. IEEE, Dec. 2014, pp. 220–223.
- [21] K. Deb, “Multi-objective genetic algorithms: problem difficulties and construction of test problems.” *Evolutionary computation*, vol. 7, no. 3, pp. 205–30, Jan. 1999.
- [22] Y. Yu and V. Prasanna, “Power-aware resource allocation for independent tasks in heterogeneous real-time systems,” in *Ninth International Conference on Parallel and Distributed Systems, 2002. Proceedings.* IEEE Comput. Soc, 2002, pp. 341–348.
- [23] H. Quinn, M. Leeser, and L. Smith King, “Dynamo: a runtime partitioning system for FPGA-based HW/SW image processing systems,” *Journal of Real-Time Image Processing*, vol. 2, no. 4, pp. 179–190, 2007.

- [24] D. Llamocca and M. Pattichis, “A Dynamically Reconfigurable Pixel Processor System Based on Power/Energy-Performance-Accuracy Optimization,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 23, no. 3, pp. 488–502, Mar. 2013.
- [25] C. G. Blair and N. M. Robertson, “Event-Driven Dynamic Platform Selection for Power-Aware Real-Time Anomaly Detection in Video,” in *International Conference on Computer Vision Theory and Applications (VISAPP 2014)*, Lisbon, Jan. 2014.
- [26] N. Dalal, “Finding People in Images and Videos,” PhD Thesis, Institut National Polytechnique de Grenoble / INRIA Grenoble, 2006.
- [27] M. Everingham, L. Gool, C. K. I. Williams, J. Winn, and A. Zisserman, “The Pascal Visual Object Classes (VOC) Challenge,” *International Journal of Computer Vision*, vol. 88, no. 2, pp. 303–338, Sep. 2009.
- [28] C. Stauffer and W. Grimson, “Adaptive background mixture models for real-time tracking,” in *Proceedings of the 1999 IEEE Conference on Computer Vision and Pattern Recognition*, vol. 99, 1999, pp. 246–252.
- [29] Z. Zivkovic, “Improved adaptive Gaussian mixture model for background subtraction,” in *Proceedings of the 17th International Conference on Pattern Recognition, 2004.*, vol. 2. IEEE, 2004, pp. 28–31.
- [30] Home Office Centre for Applied Science and Technology, *Imagery Library for Intelligent Detection Systems (I-LIDS) User Guide*, 4th ed. UK Home Office, 2011.

Calum G. Blair received the M.Eng. degree from Glasgow University, U.K. in 2009 and the Eng.D. degree from Heriot-Watt University, U.K., in 2014. He is currently a Postdoctoral Research Fellow with the Institute for Digital Communications in the School of Engineering at University of Edinburgh. His research interests include real-time implementation of image and signal processing algorithms on parallel architectures, such as multicore, GPU and FPGA, and



methods for obtaining high performance from embedded systems while minimising power consumption.

Neil M. Robertson (SM10) received the M.Sci. degree from Glasgow University, Glasgow, U.K., in 2000, and the D.Phil. degree from Oxford University, Oxford, U.K., in 2006.

He is principal investigator of the Visionlab at Heriot-Watt University and an Honorary Fellow of the School of Engineering, University of Edinburgh. He co-leads the EPSRC Edinburgh Centre for Robotics and the EPSRC/DSTL University Defence Research Centre. His work centres on human behavior recognition, computer vision and multi-modal sensor fusion with his research team. From 2000 to 2007, he worked in the U.K. Scientific Civil Service with DERA, and then QinetiQ. He held a prestigious 1851 Royal Commission Fellowship at Oxford University (2003-2006) in the Robotics Research Group.